

Online Research @ Cardiff

This is an Open Access document downloaded from ORCA, Cardiff University's institutional repository: <https://orca.cardiff.ac.uk/id/eprint/102045/>

This is the author's version of a work that was submitted to / accepted for publication.

Citation for final published version:

Huang, Haozhi, Fang, Xiaonan, Ye, Yufei, Zhang, Songhai and Rosin, Paul L. ORCID: <https://orcid.org/0000-0002-4965-3884> 2017. Practical automatic background substitution for live video. Computational Visual Media 3 (3) , pp. 273-284. 10.1007/s41095-016-0074-0 file

Publishers page: <http://dx.doi.org/10.1007/s41095-016-0074-0>
<<http://dx.doi.org/10.1007/s41095-016-0074-0>>

Please note:

Changes made as a result of publishing processes such as copy-editing, formatting and page numbers may not be reflected in this version. For the definitive version of this publication, please refer to the published source. You are advised to consult the publisher's version if you wish to cite this paper.

This version is being made available in accordance with publisher policies.

See

<http://orca.cf.ac.uk/policies.html> for usage policies. Copyright and moral rights for publications made available in ORCA are retained by the copyright holders.



Practical automatic background substitution for live video

Haozhi Huang¹, Xiaonan Fang¹, Yufei Ye¹, Songhai Zhang¹ (✉), and Paul L. Rosin²

© The Author(s) 2016. This article is published with open access at Springerlink.com

Abstract In this paper we present a novel automatic background substitution approach for live video. The objective of background substitution is to extract the foreground from the input video and then combine it with a new background. In this paper, we use a color line model to improve the Gaussian mixture model in the background cut method to obtain a binary foreground segmentation result that is less sensitive to brightness differences. Based on the high quality binary segmentation results, we can automatically create a reliable trimap for alpha matting to refine the segmentation boundary. To make the composition result more realistic, an automatic foreground color adjustment step is added to make the foreground look consistent with the new background. Compared to previous approaches, our method can produce higher quality binary segmentation results, and to the best of our knowledge, this is the first time such an automatic and integrated background substitution system has been proposed which can run in real time, which makes it practical for everyday applications.

Keywords background substitution; background replacement; background subtraction; alpha matting

1 Introduction

Background substitution is a fundamental post-processing technique for image and video editing. It

has extensive applications in video composition [1, 2], video conferencing [3, 4], and augmented reality [5]. The process of background substitution can be basically separated into two steps. The first step is to extract the foreground from the input video, and the second step is to combine the original foreground with the new background. Given limited computational resources and time, it is even more challenging to achieve satisfactory background substitution results in real time for live video. In this paper, we focus on background substitution for live video and especially live chat video, in which the camera is monocular and static, and the background is also basically static.

Foreground segmentation, also known as matting, is a fundamental problem. Formally, foreground segmentation takes as input an image I , which is assumed to be a composite of a foreground image F and a background image B . The color of the i th pixel can be represented as a linear combination of the foreground and background colors, where α represents the opacity value:

$$I_i = \alpha_i F_i + (1 - \alpha_i) B_i \quad (1)$$

This is an ill-posed problem which needs assumptions or extra constraints to become solvable.

Generally, existing work on foreground segmentation can be categorized into automatic approaches or interactive approaches. Automatic approaches usually assume that the camera and background are static, and a pre-captured background image is available. They try to model the background using either generative methods [6–9], or non-parametric methods [10, 11]. Those pixels which are consistent with the background model are labeled as background, and the remainder are labeled as foreground. Some recent works incorporate a conditional random field to include color, contrast, and motion cues, and use graph-

1 Department of Computer Science, Tsinghua University, Beijing, 100084, China. E-mail: H. Huang, huanghz08@gmail.com; X. Fang, wwjpromise@163.com; Y. Ye, yeyf13.judy@gmail.com; S. Zhang, shz@tsinghua.edu.cn (✉).

2 School of Computer Science and Informatics, Cardiff University, Cardiff, CF24 3AA, UK. E-mail: rosinpl@cardiff.ac.uk.

Manuscript received: 2016-08-18; accepted: 2016-12-20

cut to solve an optimization problem [12–14]. Most online automatic approaches only produce a binary foreground segmentation instead of fractional opacities for the sake of time, and then use feathering [12] or border matting [13] to compute approximate fractional opacities along the boundary. Feathering is a relatively crude, but efficient, technique that fades out the foreground at a fixed rate. Border matting is an alpha matting method that is significantly simplified to only collect the nearby foreground/background samples for each unknown pixel to allow fitting of a Gaussian distribution, which is then used to estimate the alpha value for that pixel. Although border matting also uses dynamic programming to minimize an energy function that encourages alpha values varying smoothly along the boundary, the result of border matting is far from globally optimal. On the other hand, interactive approaches have been proposed to handle more complicated camera motion [1, 15, 16]. Since strictly real-time performance is unnecessary for such applications, they compute more precise fractional opacities along the segmentation boundary from the beginning. Such methods require the user to draw some strokes or a trimap in a few frames to indicate whether a pixel belongs to the foreground/background/unknown region. They then solve for the alpha values in the unknown region and propagate the alpha mask to other frames.

In contrast to the large amount of foreground segmentation publications, there are fewer studies on techniques for compositing the original foreground and a new background for background substitution. Since the light sources of the original video and the new background may be drastically different, directly copying the foreground to the new background will not achieve satisfactory results. Some seamless image composition techniques [17, 18] may seem relevant at first glance, but they require the original and new backgrounds to be similar. Other color correction techniques based on color constancy [19–22] are more suitable in our context. Color constancy methods first estimate the light source color of the image, and then adjust pixel colors according to the specified hypothetical light source color.

In this paper, we present a novel practical automatic background substitution system for live

video, especially live chat video. Since real-time performance is necessary and interaction is inappropriate during live chat, our method is designed to be efficient and automatic. We first accomplish binary foreground segmentation by a novel method which is based on *background cut* [12]. To make the segmentation result less sensitive to brightness differences, we introduce a simplified version of the *color line model* [23] during the background modeling stage. Specifically, we build a color line for each background pixel and allow larger variance along the color line than in the perpendicular direction. We also include a more recent promising alpha matting method [24] to refine the segmentation boundary instead of feathering [12] or border matting [13]. To maintain real-time performance when including such a complicated alpha matting process, we perform foreground segmentation at a coarser level and then use simple but effective bilinear upsampling to generate a foreground mask for the finer level. After foreground segmentation, in order to compensate for any lighting difference between the input video and the new background, we estimate the color of the light sources in both the input video and new background, and then adjust the foreground color based on the color ratio of the light sources. This color compensation process follows the same idea as the white-patch algorithm [25], but to our knowledge this is the first time such color compensation has been applied to background substitution. Compared to previous approaches, thanks to its invariance to luminance changes, the binary segmentation result of our method is more accurate, and, thanks to the alpha matting border refinement and foreground color compensation, the appearance of the foreground in our result is more compatible with the new background.

In summary, the main contributions of our paper are:

- A novel practical automatic background substitution system for live video.
- Introduction of a color line model in conjunction with a Gaussian mixture model in the background modeling stage, which makes the foreground segmentation result less sensitive to brightness differences.
- Application of a color compensation step to

background substitution, which makes the inserted foreground look more natural in the new background.

2 Related work

2.1 Automatic video matting

Unlike interactive video matting methods [1, 15, 16, 26], which need user interaction during video playback, automatic video matting is more appropriate for live video. The earliest kind of automatic video matting problem is constant color matting [27], which uses a constant backing color, often blue, so is usually called *blue screen matting*. Although excellent segmentation results can be achieved by blue screen matting, it needs extra equipment including a blue screen and careful setting of light sources. More recent video matting methods loosen the requirement for the background to have constant color, and only assume that the background can be pre-captured and remains static or only contains slight movements. They model the background using either generative methods, such as a Bayesian model [6], a self-organized map [7], a Gaussian mixture model [8], independent component analysis [9], a foreground–background mixture model [28], or non-parametric methods [10, 11]. Such models allow prediction of the probability of a pixel belonging to the background. These methods can create holes in the foreground and noise in the background if the colors of the foreground and background are similar, because they only make local decisions. Some recent techniques utilize the power of graph-cut to solve an optimization problem based on a conditional random field using color, contrast, and motion cues [12–14]; they are able to create more complete foreground masks since they constrain the alpha matte to follow the original image gradient. Other work [29] focuses on foreground segmentation for animation. In our case, in order to acquire real-time online matting for live video, it is inappropriate to include motion cues. Thus our model is only based on color and contrast, like the work of Sun et al. [12]. We also find that stronger shadow resistance can be achieved by employing a *color line model* [23]. Another drawback of existing online methods is that they only acquire a binary foreground segmentation and

then use approximate border refinement techniques such as feathering [12] or border matting [13] to compute fractional opacities along the boundary. In this paper, we will show that a more precise alpha matting technique can be incorporated while real-time performance can still be achieved by performing foreground segmentation at a coarser level and then using simple bilinear upsampling to generate a finer level foreground mask.

2.2 Interactive video matting

Interactive video matting is another popular video matting approach. It no longer requires a known background and static camera, and takes a user drawn trimap or strokes to tell if a pixel belongs to the foreground/background/unknown region. For images, previous methods are often sampling-based [30], affinity-based [24], or a combination of both [31], computing alpha values for the unknown region based on the known region information. For video, Chuang et al. [15] use optical flow to propagate the trimap from one frame to another. *Video SnapCut* [1] maintains a collection of local classifiers around the object boundary. Each classifier subsequently solves a local binary segmentation problem, and classifiers of one frame are propagated to subsequent frames according to motion vectors estimated between frames. However, they need to take all frames all at once to compute reliable motion vectors, which takes a huge amount of time, so are unsuitable for online video matting. Gong et al. [16] use two competing one-class support vector machines (SVMs) to model the background and foreground separately for each frame at every pixel location, use the probability values predicted by the SVMs to estimate the alpha matte; they update the SVMs over time. Near real-time performance is available with the help of a GPU, but they still need an input user trimap and an extra training stage, so are inconvenient for live video applications.

There are three main categories of methods for color adjustment to improve the realism of image composites. The first category focuses on color consistency or color harmony. For example, Wong et al. [32] adjust foreground colors to be consistent with nearby surrounding background pixels, but their method fails when nearby background pixels do not correctly represent the overall lighting conditions.

Cohen-Or et al. [33] and Kuang et al. [34] consider overall color harmony based on either aesthetic rules or models learned from a dataset, but they tend to focus on creating aesthetic images rather than realistic images. The second category of methods focuses on seamless cloning based on solving a Poisson equation or coordinate interpolation [2, 17, 18, 35, 36]. A major assumption in these approaches is that the original background is similar to the new background, which we cannot guarantee in our application. The third category of methods is based on color constancy, estimating the illumination of the image first and then adjusting colors accordingly [19–22]. In this paper, we utilize the most basic and popular color constancy method, the *white-patch algorithm* [25], to estimate the light source color, since we need its efficiency for real-time application.

3 Our approach

3.1 Overview

We now outline our method. The pipeline can be separated into three steps: foreground segmentation, border refinement, and final composition. Firstly, for the foreground segmentation step, we suppose the background can be pre-captured and maintains static. Inspired by *background cut* [12], we build a global Gaussian mixture background model, local single Gaussian background models at all pixel locations, and a global Gaussian mixture foreground model. But unlike background cut, instead of using an isotropic variance for the local single Gaussian background models, we make the variance along the *color line* larger than that in the direction perpendicular to the color line. Here the concept of a color line is borrowed from Ref. [23]. The original color line model built multiple curves to represent all colors of the whole image, and assumed that colors from the same object lie on the same curve. To check which curve a pixel belongs to is a time consuming process. In order to achieve real-time performance, we adapt the color line model to a much simpler and more efficient version. In our basic version of the color line model, for each pixel we build a single curve color line model, which avoids the process of matching a pixel to one of the curves in the multiple curve model. Furthermore, instead of fitting a curve, we fit a straight line

that intersects the origin in RGB space, which means we ignore the non-linear transform of the camera sensor. Our experiments show this simplified model to be sufficient and effective. By utilizing this color line model, we can avoid misclassifying background pixels which undergo color changes due to a shadow passing by, since color changes caused by shadows still remain along the color line. Using this background cut model, we can build an energy function that can be optimized by graph-cut to give a binary foreground segmentation matte. Secondly, we carry out border refinement for this binary foreground matte. Specifically, we use morphological operations to mark the border pixels between foreground and background. Considering these border pixels to be the unknown region results in a trimap. We then carry out closed-form alpha matting [24], which computes fractional alpha values for these border pixels. It is important to emphasize that, only when the binary foreground segmentation result is essentially correct, can a trimap be automatically computed reliably in this way. Lastly, to perform composition, we estimate the light source colors of the original input video and the new background separately, and adjust the foreground colors accordingly to make the foreground look more consistent with the new background.

3.2 Foreground segmentation

3.2.1 Basic background cut model

In this section we briefly describe the background cut model proposed in Ref. [12]. The background cut algorithm takes a video and a pre-captured background as input, and the output is a sequence of binary foreground masks, in which each pixel r is labeled 0 if it belongs to the background or 1 otherwise. Background cut solves the foreground segmentation problem frame by frame. For each frame, the process of labeling can be transformed into solving a global optimization problem. The energy function to be minimized is in the form of a *conditional random field*:

$$E(X) = \sum_r E_d(x_r) + \lambda_1 \sum_{r,s} E_c(x_r, x_s) \quad (2)$$

where $X = \{x_r\}$, x_r denotes the label value, r, s are neighbouring pixels in one frame, E_d (the data term) represents per-pixel energy, and E_c is a contrast term computed from neighbouring pixels. Here λ_1 is a predefined constant balancing E_d and E_c , which is

empirically set to 30 in our experiments. This is a classical energy function which can be minimized by graph-cut [37].

Now we explain how to construct E_d and E_c . First we model the foreground and the background using Gaussian models. For the foreground, we build a global Gaussian mixture model (GMM). For the background, we not only build a global GMM, but also a local single Gaussian distribution model at each pixel location (a per-pixel model). The two global GMMs are defined as

$$p(v_r|x_r = i) = \sum_{k=1}^{k_i} w_k^i N(v_r|\mu_k^i, \Sigma_k^i), \quad i = 0, 1 \quad (3)$$

where $i = 0$ and $i = 1$ stand for background and foreground respectively, v_r denotes the color of pixel r , k_i denotes the number of mixture components, w_k^i denotes the weight of the k th component, N denotes the Gaussian distribution, μ_k^i denotes the mean, and Σ_k^i denotes the covariance matrix. The single Gaussian distribution at every pixel location is defined as

$$p_s(v_r) = N(v_r|\mu_r^s, \Sigma_r^s) \quad (4)$$

where $\Sigma_r^s = \sigma_r^s I$, so, following Ref. [12], the variance of the per-pixel model is isotropic. The background global GMM and the background per-pixel model are initialized using pre-captured background data. The foreground global GMM is initialized using pixels whose probabilities are lower than a threshold in the background model. After initialization, these Gaussian models are updated frame by frame according to the segmentation results.

Based on the Gaussian models, the data term E_d is defined as

$$E_d(x_r) = \begin{cases} -\log(\lambda_2 p(v_r|x_r) + (1-\lambda_2)p_s(v_r)), & x_r = 0 \\ -\log p(v_r|x_r), & x_r = 1 \end{cases} \quad (5)$$

Here λ_2 is a predefined constant balancing the global GMM and the local per-pixel model, which is empirically set to 0.1 in our experiments. The contrast term is

$$E_c(x_r, x_s) = |x_r - x_s| \exp(-\beta \|v_r - v_s\|^2 / d_B(r, s)) \quad (6)$$

$$d_B(r, s) = 1 + (\|v_r^B - v_s^B\|/K)^2 \exp(-z_{rs}^2/\sigma_z) \quad (7)$$

where $d_B(r, s)$ is a contrast attenuation term proportional to the contrast with respect to the background, $z_{rs} = \max(\|v_r - v_r^B\|, \|v_s - v_s^B\|)$ measures the dissimilarity between the pre-captured

background and the current frame, and β , K , and σ_z are predefined constants. In our experiments, we set $\beta = 0.005$, $K = 1$, $\sigma_z = 10$. The introduction of the contrast attenuation term causes E_c to rely on the contrast from the foreground instead of the background.

The energy function in Eq. (2) can be optimized using the graph-cut algorithm [37]. For more details of the model, please refer to Ref. [12]. One major drawback of this background cut model is that, when the color of a background pixel changes due to changes in illumination, it will have extremely low probability in the per-pixel model, which will cause the pixel to be misclassified as foreground instead of background.

3.2.2 Background cut with color line model

Now we explain how the color line model [23] can improve the effectiveness of the background cut model in the presence of shadows.

Starting from the basic color line model, we make the assumption that colors of a certain material under different intensities of light form a linear color cluster that intersects the origin in RGB space. Suppose the average color at a pixel location is $\mu_r^s = (r, g, b)$. When the illumination of the same pixel location changes, its color will also change from μ_r^s to v_r . According to the color line model, v_r will approximately lie on the line connecting the origin and μ_r^s in the RGB color space. With this insight, we can decompose v_r as

$$v_r = v_\perp + v_\parallel \quad (8)$$

such that $v_\perp \perp \mu_r^s$ and $v_\parallel \parallel \mu_r^s$. Define:

$$f(v_r, \mu_r^s) = N(\|v_\perp\| \mid 0, \sigma_{pe}) N(\|v_\parallel\| \mid \|\mu_r^s\|, \sigma_{pa}) \quad (9)$$

where σ_{pe} and σ_{pa} are the respective variances of the Gaussian distributions for the perpendicular direction and parallel direction. Then the per-pixel single Gaussian distribution Eq. (4) is modified to be

$$p_s(v_r) = f(v_r, \mu_r^s) \quad (10)$$

As discussed before, the color of an object is more likely to fluctuate in the parallel direction than in the perpendicular direction. Therefore, we set $\sigma_{pe} = \sigma_r^s$, $\sigma_{pa} = \lambda_3 \sigma_{pe}$, $\lambda_3 > 1$ to constrain variance in the perpendicular direction and tolerate variance in the parallel direction, which gives our model a strong resistance to shadow. Here we do not build a global color line model as in Ref. [23], which uses multiple color lines for the whole image to replace the global

GMM, because it takes a long time to determine which line each pixel belongs to when the number of lines is large (e.g., a model with 40 lines is used in Ref. [23]), precluding real-time performance.

3.3 Border refinement

After graph-cut, we add an extra hole filling step by applying the morphological close operation to fill small holes in the foreground mask. See Fig. 1 for an example. However, we currently still have a binary foreground matte (see Fig. 1(d)). In this subsection, we explain how to automatically compute fractional alpha values for the segmentation border.

First, we automatically generate a mask covering the segmentation border as the unknown region:

$$U_i = 1 - (\text{erode}(F)_i \wedge \text{erode}(B)_i) \quad (11)$$

Here U_i denotes the value of the i th pixel of the unknown mask, $\text{erode}()$ denotes the morphological erosion operation, F is the binary foreground matte, and B is the binary background matte where $B_i = 1 - F_i$. The morphological operation radius is set to 2 for 640×480 input. The eroded foreground mask, eroded background mask, and unknown region mask are separately painted in white, black, and gray in the final trimap. Using this trimap with one of the most popular alpha matting methods [24], we calculate the fractional alpha values for the unknown region. See Fig. 2 for an example of the generated trimap and alpha matting result.

3.4 Composition

For an ideal final composition, the new composite image should be

$$I_{\text{new}} = \alpha F_{\text{old}} + (1 - \alpha) B_{\text{new}} \quad (12)$$

Here I_{new} denotes the new composite image, F_{old}



Fig. 2 (a) Automatically generated trimap. (b) Alpha matting result.

denotes the original foreground, and B_{new} denotes the new background (Fig. 3(c)). For previous methods whose pre-captured background (Fig. 3(a)) is unavailable, F_{old} is approximated by I_{old} :

$$I_{\text{new}} = \alpha I_{\text{old}} + (1 - \alpha) B_{\text{new}} \quad (13)$$

However, in our case, since the pre-captured background B_{old} is available, we can calculate the original foreground more accurately:

$$F_{\text{old}} = [I_{\text{old}} - (1 - \alpha) B_{\text{old}}] / \alpha \quad (14)$$

This gives a final composition formula:

$$I_{\text{new}} = I_{\text{old}} + (1 - \alpha)(B_{\text{new}} - B_{\text{old}}) \quad (15)$$

Directly applying the above composition will create unrealistic results due to the difference in light source colors between the original input and the new background. Thus, we propose a color compensation process to deal with this problem.

First, we need to estimate the light source colors of the original input video and the new background image. The *white-patch method* [25], a popular color constancy method, assumes that the highest values in each color channel represent the presence of white in the image. In this paper, we use the variant of the white-patch method designed for CIE-Lab space, a color space that is naturally designed to separate lightness and chroma. We first calculate

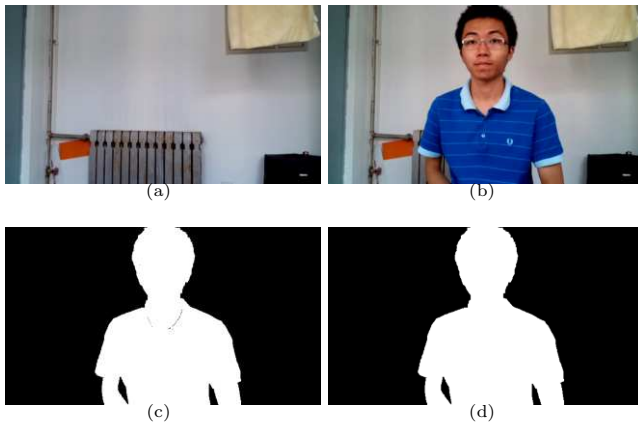


Fig. 1 (a) Pre-captured background. (b) One frame of the input video. (c) Binary foreground matte after graph-cut. (d) Foreground matte after filling holes.

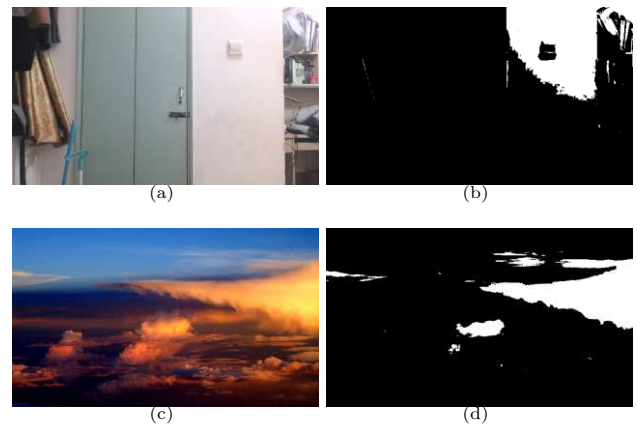


Fig. 3 (a) Pre-captured background. (b) Estimated light source mask of the pre-captured background (a). (c) New background. (d) Estimated light source mask of the new background (c).

the accumulated histogram in the lightness channel L of an image in CIE-Lab space, and consider those 10% pixels with the largest lightness values to be the white pixels. Figure 3 shows an example of the light source masks. The estimated light source color is then computed as the mean color value of all light source pixels. Denote the estimated light source color of the input video as c_{old} , that of the new background image as c_{new} . Then the new composite image after color compensation is

$$I_{new} = rI_{old} + (1 - \alpha)(B_{new} - rB_{old}) \quad (16)$$

$$r = c_{new}/c_{old} \quad (17)$$

Figure 4 compares results with and without light source color compensation. We can clearly see that the result with color compensation is more realistic.

4 Results and discussion

In this section, we report results generated under different conditions. All results shown were generated using fixed parameters.

Results for different frames of the same input video. Figure 5 shows that our method can create generally good background substitution results for different frames, no matter what the gesture is. Sometimes there may be residual background between the fingers (e.g., Fig. 5(c)) due

to the hole-filling post-processing, but it does not do much harm to the overall effect.

Results for different input videos. Figure 6 shows that our method can deal with different kinds of foreground and background. Color compensation works fine for various lighting condition. Although the matting border is not 100% perfect for Fig. 6(b) due to confusion of hair and background, the composition result is generally good.

Comparison with previous methods. We compare various methods: *fuzzy Gaussian* [38], *adaptive-SOM* [7], *background cut* [12] using RGB color space and CIE-Lab color space, and our color line model. To implement the fuzzy Gaussian and adaptive-SOM methods, we used the code in the BGSLibrary [39]. There are also other background subtraction methods in the BGSLibrary, we choose these two methods because they show the most promising results under real-time conditions. Figure 7 shows foreground masks created by different methods. After the person walks into the picture, some shadow will be cast onto the wall. The fuzzy Gaussian and adaptive-SOM methods create a lot of noise and holes since they do not utilize gradient information between neighbouring pixels. Background cut used in RGB color space does a better job by using the graph-cut model to introduce gradient information. However, it is sensitive to brightness differences, which causes shadow to be misclassified as foreground. If we set the variance of the Gaussian to be larger to tolerate some shadow, part of the true foreground is then misclassified as background. Background cut in CIE-Lab color space also suffers from the same issue. Although allowing a larger variance in the L channel can give

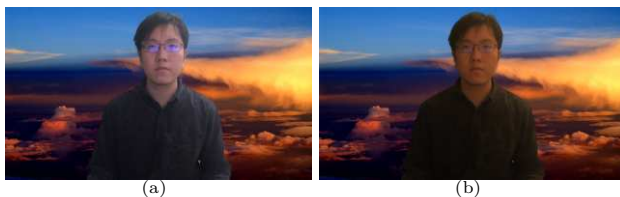


Fig. 4 (a) Composite result without color compensation. (b) Composite result with color compensation.

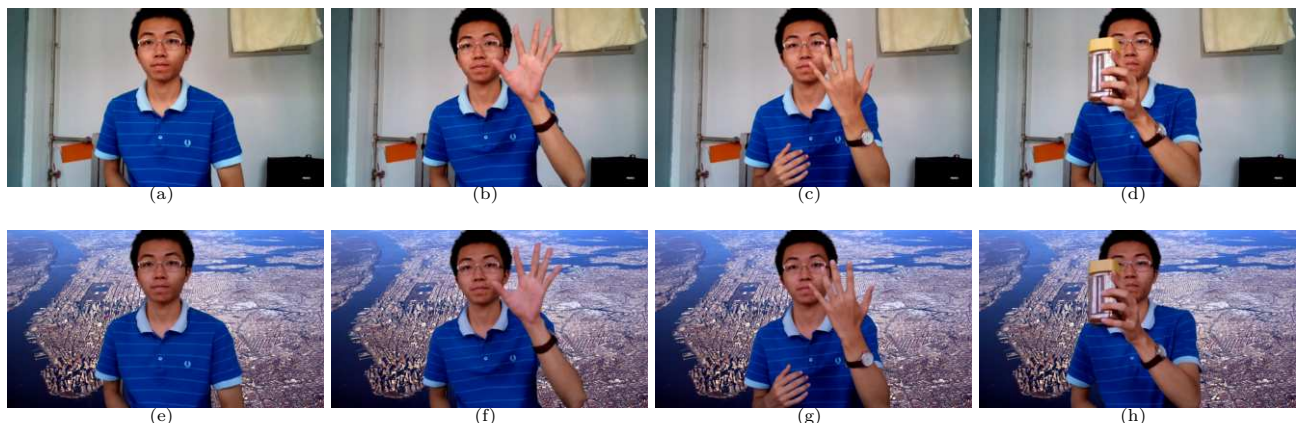


Fig. 5 (a)–(d) Input video frames. (e)–(h) Background substitution results.

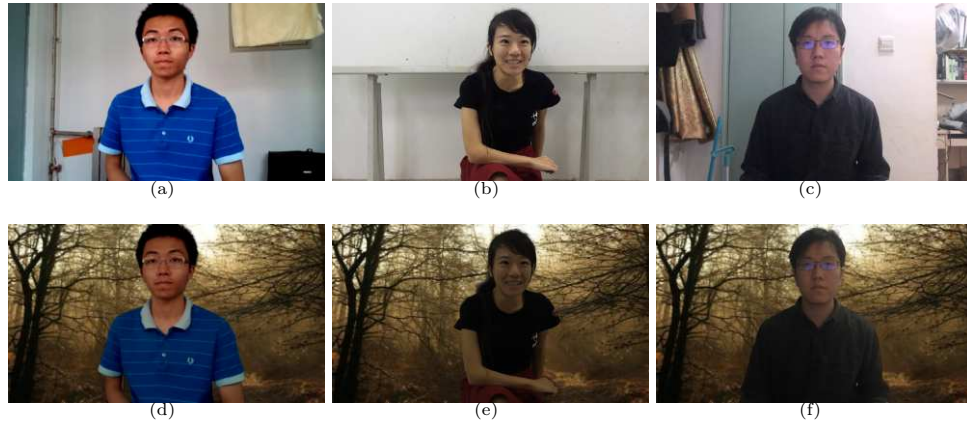


Fig. 6 (a)–(c) Input frames from different videos. (d)–(f) Background substitution results.

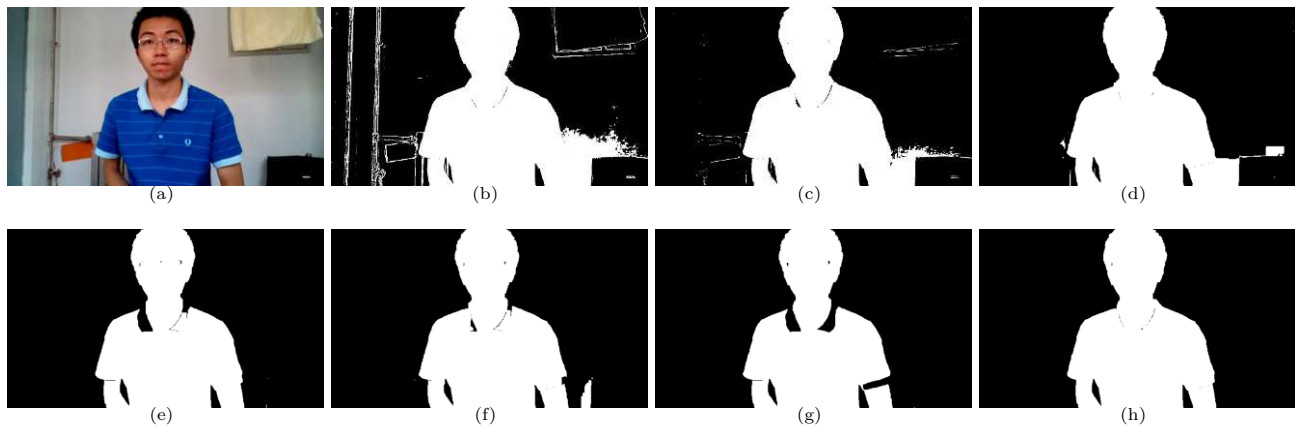


Fig. 7 (a) Input frame. (b) Foreground mask created by fuzzy Gaussian. (c) Adaptive-SOM result. (d) Background cut in RGB space result with a low variance ($\sigma_r^s = (5/255)^2$) Gaussian model. (e) Background cut in RGB space result with a larger variance ($\sigma_r^s = (20/255)^2$). (f) Background cut in CIE-Lab space result with $\sigma_L = \sigma_a = \sigma_b = (5/255)^2$. (g) Background cut in CIE-Lab space result with larger variance in the L channel ($\sigma_L = 5 * (5/255)^2$). (h) Result of our method ($\sigma_{pe} = (10/255)^2, \sigma_{pa} = 10 * (10/255)^2$).

greater tolerance to brightness changes, in actual test cases, even when we only increase the variance in the L channel by a small amount, part of the collar disappears. In contrast, using our color line model with background cut constantly creates a better foreground segmentation result.

To further quantitatively evaluate the comparison, we created a large number of “ground truth” foreground masks following a similar approach to one in Ref. [40]. The key idea is to use some balls as the moving foreground objects, and use a circle detection technique to detect the balls, which will automatically create “ground truth” masks for evaluation of our foreground segmentation methods. Specifically, we first calculate the difference image between the pre-captured background and the current frame (where one or more balls appear). Then we perform circle detection using the Hough transform [41] on the difference image, which generally produces reliable and accurate detection

results. Finally, we manually eliminate the small number of outliers that occur when circle detection fails. In total, 4105 frames and their circle detection results are collected as the ground truth. Figure 8 shows a few examples. We did not use the ground truth from the *VideoMatting* benchmark [42], because their synthetic test images do not have shadows on the background, which is one of the fundamental aspects we wish to test. Using the generated ground truth, we tested different methods including fuzzy Gaussian, adaptive-SOM, background cut using RGB color space and CIE-Lab color space, and our color line model. For fuzzy Gaussian and adaptive-SOM, we used the default parameters provided by the BGSLibrary. For the background cut method, we tested several parameters and gave results with the highest F1 score. Table 1 shows that background cut with our color line model achieves the highest F1 score, the CIE-Lab space method follows closely, and others are

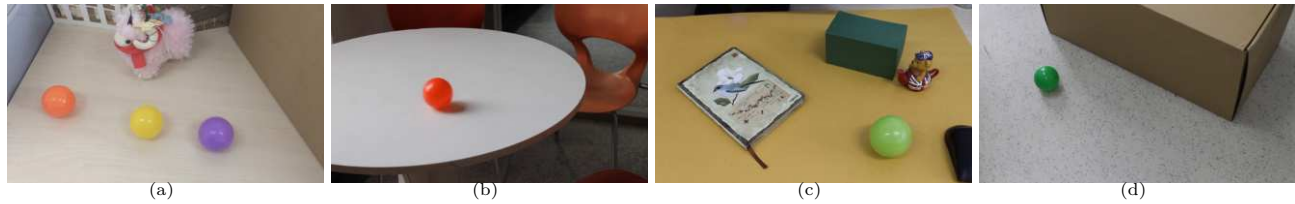


Fig. 8 Example frames for creating ground truth.

Table 1 Method comparison on ground truth dataset

| Method | Precision | Recall | F1 |
|----------------|-----------|--------|-------|
| Fuzzy Gaussian | 0.252 | 0.993 | 0.402 |
| Adaptive-SOM | 0.510 | 0.963 | 0.667 |
| BC-RGB | 0.839 | 0.962 | 0.896 |
| BC-Lab | 0.900 | 0.968 | 0.933 |
| BC-Colorline | 0.907 | 0.964 | 0.935 |

substantially worse. However, as we have already shown in Fig. 7, CIE-Lab space has an obvious drawback in actual application scenarios. We also tested an outdoor scene with different methods to show the effectiveness of our model: see Fig. 9. In conclusion, our color line model generally creates a better foreground segmentation boundary, and is effective at coping with differences in brightness.

Results with new background. We also tested our color compensation method using new backgrounds with different light sources. In Fig. 10, the first row shows the new input backgrounds, and the second row shows the light source pixel masks. The third row contains the composition results; we can see that the color of the foreground varies correctly according to different backgrounds.

Acceleration. Although we restrict the alpha matting computation to a very small region, it is still computationally expensive. In order to enable our algorithm to run in real time, we first downsample the input frames by a scale of two, carry out foreground cut and alpha matting on the downsampled images, and then upsample the matting result to the original scale. We finish the final composition step at the original scale. We call this process “sampling acceleration”. As we can see in Fig. 11, the matting result using sampling acceleration is very similar to the result produced by processing the full frames. If we do not use alpha matting to refine the border, the border is jagged (see Fig. 11(c)).

Performance. We have implemented our method in C++ on a PC with an Intel 3.4 GHz Core i7-3770 CPU. For a 640×480 input video, our background substitution program can run at 10 frames per second using just the CPU, and it can run at a real-time frame rate with GPU parallelization.

5 Conclusions

In this paper, we have presented a novel background

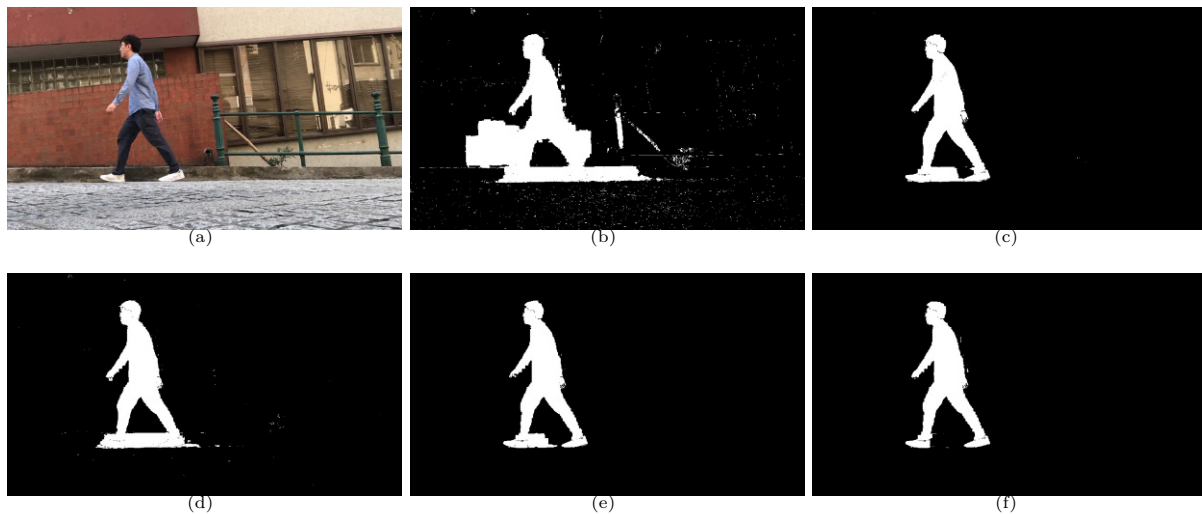


Fig. 9 (a) Input frame. (b) Background cut in RGB space result with a low variance ($\sigma_r^s = (5/255)^2$) Gaussian model. (c) Background cut in RGB space result with a larger variance ($\sigma_r^s = (20/255)^2$). (d) Background cut in CIE-Lab space result with $\sigma_L = \sigma_a = \sigma_b = (5/255)^2$. (e) Background cut in CIE-Lab space result with larger variance in the L channel ($\sigma_L = 5 * (5/255)^2$). (f) Result of our method ($\sigma_{pe} = (10/255)^2, \sigma_{pa} = 10 * (10/255)^2$).

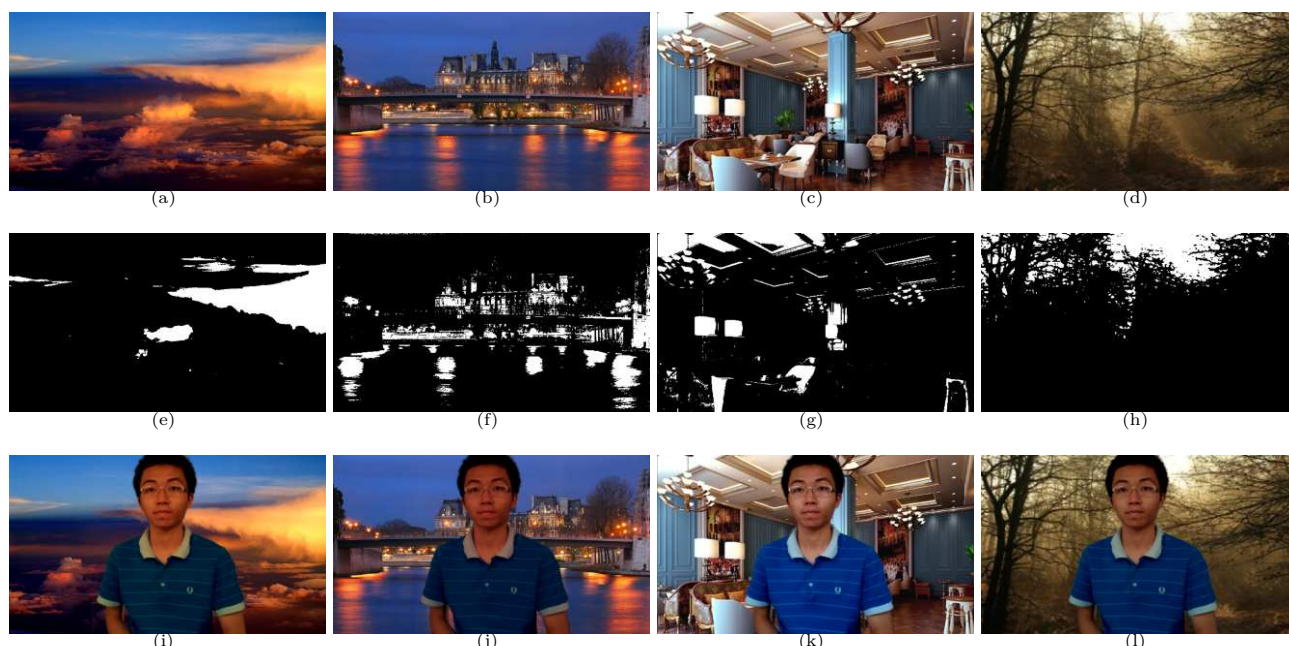


Fig. 10 (a)–(d) Input new backgrounds. (e)–(h) Estimated light source pixel mask. (i)–(l) Background substitution results.

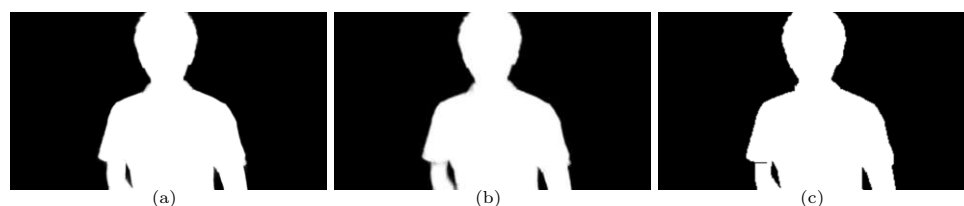


Fig. 11 (a) Matting result without sampling acceleration. (b) Matting result with sampling acceleration. (c) Foreground segmentation result with sampling acceleration but without alpha matting border refinement.

substitution method for live video. It optimizes a cost function based on Gaussian mixture models and a conditional random field, using graph-cut. A color line model is used when computing the Gaussian mixture model to make the model less sensitive to brightness differences. Before final composition, we use alpha matting to refine the segmentation border. Light source colors of the input video and new background are estimated by a simple method, and we adjust the foreground colors accordingly to give more realistic composition results. Compared to previous methods, our approach can automatically produce more accurate foreground segmentation masks and more realistic composition results, while still maintaining real-time performance.

Acknowledgements

We thank the reviewers for their valuable comments. This work was supported by the National High-Tech R&D Program of China (Project No.

2012AA011903), the National Natural Science Foundation of China (Project No. 61373069), the Research Grant of Beijing Higher Institution Engineering Research Center, and Tsinghua–Tencent Joint Laboratory for Internet Innovation Technology.

Electronic Supplementary Material Supplementary material is available in the online version of this article at <http://dx.doi.org/10.1007/s41095-016-0074-0>.

References

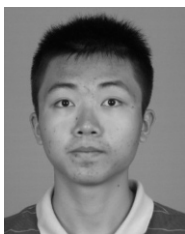
- [1] Bai, X.; Wang, J.; Simons, D.; Sapiro, G. Video SnapCut: Robust video object cutout using localized classifiers. *ACM Transactions on Graphics* Vol. 28, No. 3, Article No. 70, 2009.
- [2] Chen, T.; Zhu, J.-Y.; Shamir, A.; Hu, S.-M. Motion-aware gradient domain video composition. *IEEE Transactions on Image Processing* Vol. 22, No. 7, 2532–2544, 2013.
- [3] Liu, Z.; Cohen, M. Head-size equalization for better visual perception of video conferencing. In:

- Proceedings of the IEEE International Conference on Multimedia and Expo, 4, 2005.
- [4] Zhu, Z.; Martin, R. R.; Pepperell, R.; Burleigh, A. 3D modeling and motion parallax for improved videoconferencing. *Computational Visual Media* Vol. 2, No. 2, 131–142, 2016.
 - [5] Van Krevelen, D. W. F.; Poelman, R. A survey of augmented reality technologies, applications and limitations. *International Journal of Virtual Reality* Vol. 9, No. 2, 1–21, 2010.
 - [6] Apostoloff, N.; Fitzgibbon, A. Bayesian video matting using learnt image priors. In: Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition, Vol. 1, I-407–I-414, 2004.
 - [7] Bouwmans, T.; El Baf, F.; Vachon, B. Background modeling using mixture of Gaussians for foreground detection—A survey. *Recent Patents on Computer Science* Vol. 1, No. 3, 219–237, 2008.
 - [8] Maddalena, L.; Petrosino, A. A self-organizing approach to background subtraction for visual surveillance applications. *IEEE Transactions on Image Processing* Vol. 17, No. 7, 1168–1177, 2008.
 - [9] Tsai, D.-M.; Lai, S.-C. Independent component analysis-based background subtraction for indoor surveillance. *IEEE Transactions on Image Processing* Vol. 18, No. 1, 158–167, 2009.
 - [10] Barnich, O.; Van Droogenbroeck, M. ViBe: A universal background subtraction algorithm for video sequences. *IEEE Transactions on Image Processing* Vol. 20, No. 6, 1709–1724, 2011.
 - [11] Hofmann, M.; Tiefenbacher, P.; Rigoll, G. Background segmentation with feedback: The pixel-based adaptive segmenter. In: Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops, 38–43, 2012.
 - [12] Sun, J.; Zhang, W.; Tang, X.; Shum, H.-Y. Background cut. In: *Computer Vision—ECCV 2006*. Leonardis, A.; Bischof, H.; Pinz, A. Eds. Springer Berlin Heidelberg, 628–641, 2006.
 - [13] Criminisi, A.; Cross, G.; Blake, A.; Kolmogorov, V. Bilayer segmentation of live video. In: Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 53–60, 2006.
 - [14] Yin, P.; Criminisi, A.; Winn, J.; Essa, I. Bilayer segmentation of webcam videos using tree-based classifiers. *IEEE Transactions on Pattern Analysis and Machine Intelligence* Vol. 33, No. 1, 30–42, 2011.
 - [15] Chuang, Y.-Y.; Agarwala, A.; Curless, B.; Salesin, D. H.; Szeliski, R. Video matting of complex scenes. *ACM Transactions on Graphics* Vol. 21, No. 3, 243–248, 2002.
 - [16] Gong, M.; Qian, Y.; Cheng, L. Integrated foreground segmentation and boundary matting for live videos. *IEEE Transactions on Image Processing* Vol. 24, No. 4, 1356–1370, 2015.
 - [17] Pérez, P.; Gangnet, M.; Blake, A. Poisson image editing. *ACM Transactions on Graphics* Vol. 22, No. 3, 313–318, 2003.
 - [18] Jia, J.; Sun, J.; Tang, C.-K.; Shum, H.-Y. Drag-and-drop pasting. *ACM Transactions on Graphics* Vol. 25, No. 3, 631–637, 2006.
 - [19] Buchsbaum, G. A spatial processor model for object colour perception. *Journal of the Franklin Institute* Vol. 310, No. 1, 1–26, 1980.
 - [20] Finlayson, G. D.; Hordley, S. D.; Hubel, P. M. Color by correlation: A simple, unifying framework for color constancy. *IEEE Transactions on Pattern Analysis and Machine Intelligence* Vol. 23, No. 11, 1209–1221, 2001.
 - [21] Cheng, D.; Prasad, D. K.; Brown, M. S. Illuminant estimation for color constancy: Why spatial-domain methods work and the role of the color distribution. *Journal of the Optical Society of America A* Vol. 31, No. 5, 1049–1058, 2014.
 - [22] Cheng, D.; Price, B.; Cohen, S.; Brown, M. S. Beyond white: Ground truth colors for color constancy correction. In: Proceedings of the IEEE International Conference on Computer Vision, 298–306, 2015.
 - [23] Omer, I.; Werman, M. Color lines: Image specific color representation. In: Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition, Vol. 2, II-946–II-953, 2004.
 - [24] Levin, A.; Lischinski, D.; Weiss, Y. A closed-form solution to natural image matting. *IEEE Transactions on Pattern Analysis and Machine Intelligence* Vol. 30, No. 2, 228–242, 2008.
 - [25] Land, E. H.; McCann, J. J. Lightness and retinex theory. *Journal of the Optical Society of America* Vol. 61, No. 1, 1–11, 1971.
 - [26] Zhang, Y.; Tang, Y.-L.; Cheng, K.-L. Efficient video cutout by paint selection. *Journal of Computer Science and Technology* Vol. 30, No. 3, 467–477, 2015.
 - [27] Smith, A. R.; Blinn, J. F. Blue screen matting. In: Proceedings of the 23rd Annual Conference on Computer Graphics and Interactive Techniques, 259–268, 1996.
 - [28] Mumtaz, A.; Zhang, W.; Chan, A. B. Joint motion segmentation and background estimation in dynamic scenes. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 368–375, 2014.
 - [29] Zhang, L.; Huang, H.; Fu, H. EXCOL: An extract-and-complete layering approach to cartoon animation reusing. *IEEE Transactions on Visualization and Computer Graphics* Vol. 18, No. 7, 1156–1169, 2012.
 - [30] Gastal, E. S. L.; Oliveira, M. M. Shared sampling for real-time alpha matting. *Computer Graphics Forum* Vol. 29, No. 2, 575–584, 2010.
 - [31] Chen, X.; Zou, D.; Zhou, S.; Zhao, Q.; Tan, P. Image matting with local and nonlocal smooth priors. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 1902–1907, 2013.
 - [32] Wong, B.-Y.; Shih, K.-T.; Liang, C.-K.; Chen, H. H. Single image realism assessment and recoloring by color compatibility. *IEEE Transactions on Multimedia* Vol. 14, No. 3, 760–769, 2012.

- [33] Cohen-Or, D.; Sorkine, O.; Gal, R.; Leyvand, T.; Xu, Y.-Q. Color harmonization. *ACM Transactions on Graphics* Vol. 25, No. 3, 624–630, 2006.
- [34] Kuang, Z.; Lu, P.; Wang, X.; Lu, X. Learning self-adaptive color harmony model for aesthetic quality classification. In: *Proceedings of SPIE 9443, the 6th International Conference on Graphic and Image Processing*, 94431O, 2015.
- [35] Chen, T.; Cheng, M.-M.; Tan, P.; Shamir, A.; Hu, S.-M. Sketch2Photo: Internet image montage. *ACM Transactions on Graphics* Vol. 28, No. 5, Article No. 124, 2009.
- [36] Farbman, Z.; Hoffer, G.; Lipman, Y.; Cohen-Or, D.; Lischinski, D. Coordinates for instant image cloning. *ACM Transactions on Graphics* Vol. 28, No. 3, Article No. 67, 2009.
- [37] Boykov, Y.; Kolmogorov, V. An experimental comparison of min-cut/max-flow algorithms for energy minimization in vision. In: *Energy Minimization Methods in Computer Vision and Pattern Recognition*. Figueiredo, M.; Zerubia, J.; Jain, A. K. Eds. Springer Berlin Heidelberg, 359–374, 2001.
- [38] Sigari, M. H.; Mozayani, N.; Pourreza, H. R. Fuzzy running average and fuzzy background subtraction: concepts and application. *International Journal of Computer Science and Network Security* Vol. 8, No. 2, 138–143, 2008.
- [39] Sobral, A. BGSLibrary. 2016. Available at <https://github.com/andrewssobral/bgslibrary>.
- [40] Rosin, P. L.; Ioannidis, E. Evaluation of global image thresholding for change detection. *Pattern Recognition Letters* Vol. 24, No. 14, 2345–2356, 2003.
- [41] Kerbyson, D. J.; Atherton, T. J. Circle detection using Hough transform filters. In: *Proceedings of the 5th International Conference on Image Processing and its Applications*, 370–374, 1995.
- [42] Graphics and Media Lab. Videomattng benchmark. 2016. Available at <http://videomattng.com>.



Haozhi Huang is currently a Ph.D. student in the Department of Computer Science and Technology, Tsinghua University, China. He received his bachelor degree from Tsinghua University, in 2012. His research interests include image and video editing, and computer graphics.



Xiaonan Fang is currently an undergraduate student in the Department of Computer Science and Technology, Tsinghua University, China. His research interests include computational geometry, image processing, and video processing.



Yufei Ye is currently an undergraduate student in the Department of Computer Science and Technology, Tsinghua University, China. Her research interests lie in video processing, object detection, generative models, unsupervised learning, and representation learning.



Songhai Zhang received his Ph.D. degree from Tsinghua University, China, in 2007. He is currently an associate professor in the Department of Computer Science and Technology, Tsinghua University. His research interests include image and video processing, and geometric computing.



Paul L. Rosin is a professor in the School of Computer Science & Informatics, Cardiff University, UK. Previous posts include lecturer in the Department of Information Systems and Computing, Brunel University London, UK, research scientist at the Institute for Remote Sensing Applications, Joint Research Centre, Ispra, Italy, and lecturer at Curtin University of Technology, Perth, Australia. His research interests include the representation, segmentation, and grouping of curves, knowledge-based vision systems, early image representations, low level image processing, machine vision approaches to remote sensing, methods for evaluation of approximation algorithms, medical and biological image analysis, mesh processing, non-photorealistic rendering, and the analysis of shape in art and architecture.

Open Access The articles published in this journal are distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made.

Other papers from this open access journal are available free of charge from <http://www.springer.com/journal/41095>. To submit a manuscript, please go to <https://www.editorialmanager.com/cvmj>.